

# Likelihood Ratios for Mixtures: Continuous Approach

Simone Gittelsohn, Ph.D., [simone.gittelsohn@nist.gov](mailto:simone.gittelsohn@nist.gov)

Michael Coble, Ph.D., [michael.coble@nist.gov](mailto:michael.coble@nist.gov)

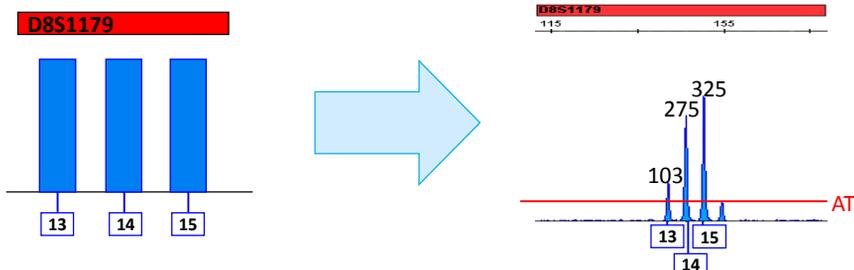
## Acknowledgement

I thank Michael Coble, Bruce Weir and John Buckleton for their helpful discussions.

## Disclaimer

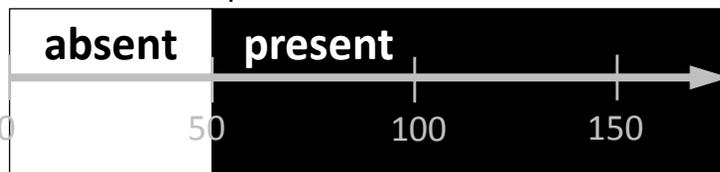
**Points of view in this presentation are mine** and do not necessarily represent the official position or policies of the National Institute of Standards and Technology.

## From Semi-continuous to Continuous



## Discrete vs. Continuous

- The observed peaks as a **discrete** variable:

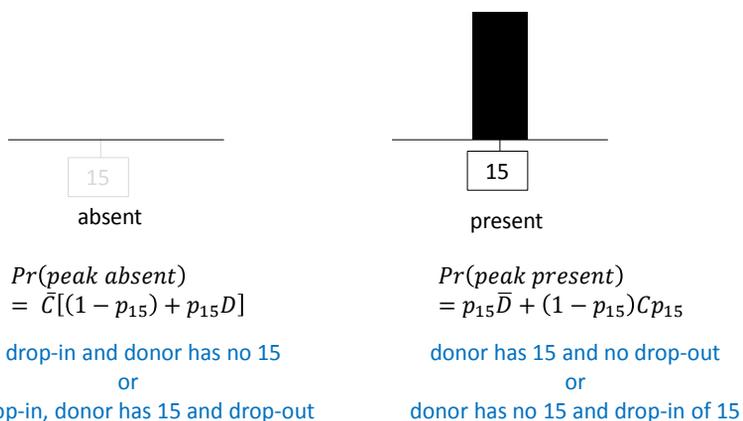


- The observed peaks as a **continuous** variable:



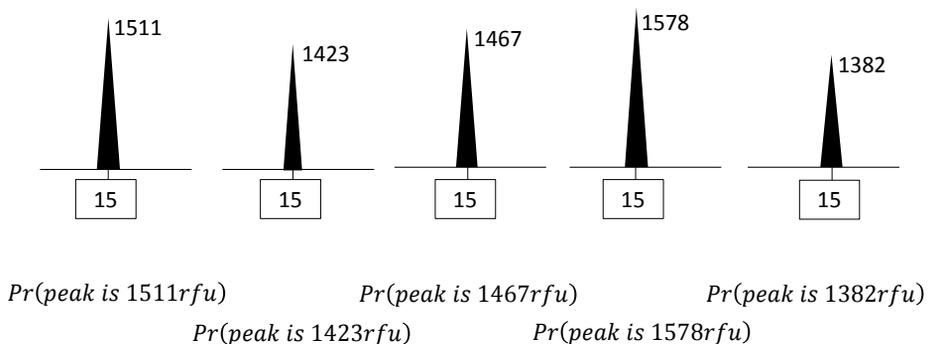
## Discrete vs. Continuous

The observed peak as a **discrete** variable:



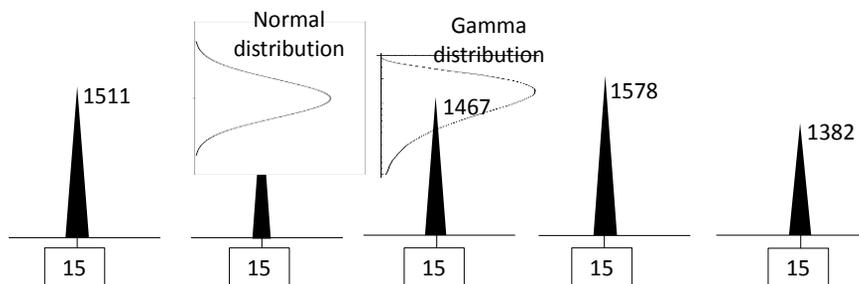
## Discrete vs. Continuous

The observed peaks as a **continuous** variable:



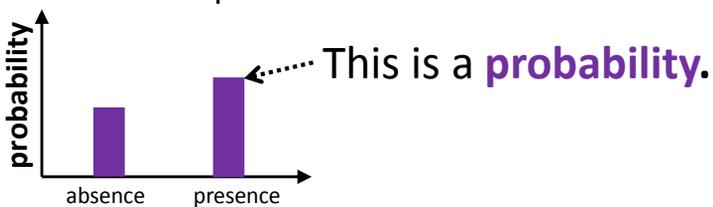
## Discrete vs. Continuous

The observed peaks as a **continuous** variable:



## Discrete vs. Continuous

- The observed peaks as a **discrete** variable:



- The observed peaks as a **continuous** variable:



Likelihood Ratio:  
from semi-continuous to continuous

$$LR = \frac{\sum_{\text{genotype sets}|H_p} \text{weight} \times \text{genotype match prob.}}{\sum_{\text{genotype sets}|H_d} \text{weight} \times \text{genotype match prob.}}$$

**stays the same**

(Hardy-Weinberg Law,  
NRC II recommendation 4.1,  
NRC II recommendation 4.2)

Likelihood Ratio:  
from semi-continuous to continuous

$$LR = \frac{\sum_{\text{genotype sets}|H_p} \text{weight} \times \text{genotype match prob.}}{\sum_{\text{genotype sets}|H_d} \text{weight} \times \text{genotype match prob.}}$$

**changes**

**semi-continuous weight:** any value between 0 and 1, including 0 and 1, assigned using probabilities of drop-out and drop-in

**continuous:** a **probability density**, generally assigned based on the results of simulations that attempt to reproduce the observed peak heights

## Likelihood Ratio

$$LR = \frac{\sum_{j=1}^m f(G_{CS}|S_j) \Pr(S_j|G_K, H_p, I)}{\sum_{j=1}^m f(G_{CS}|S_j) \Pr(S_j|G_K, H_d, I)}$$

### match probabilities

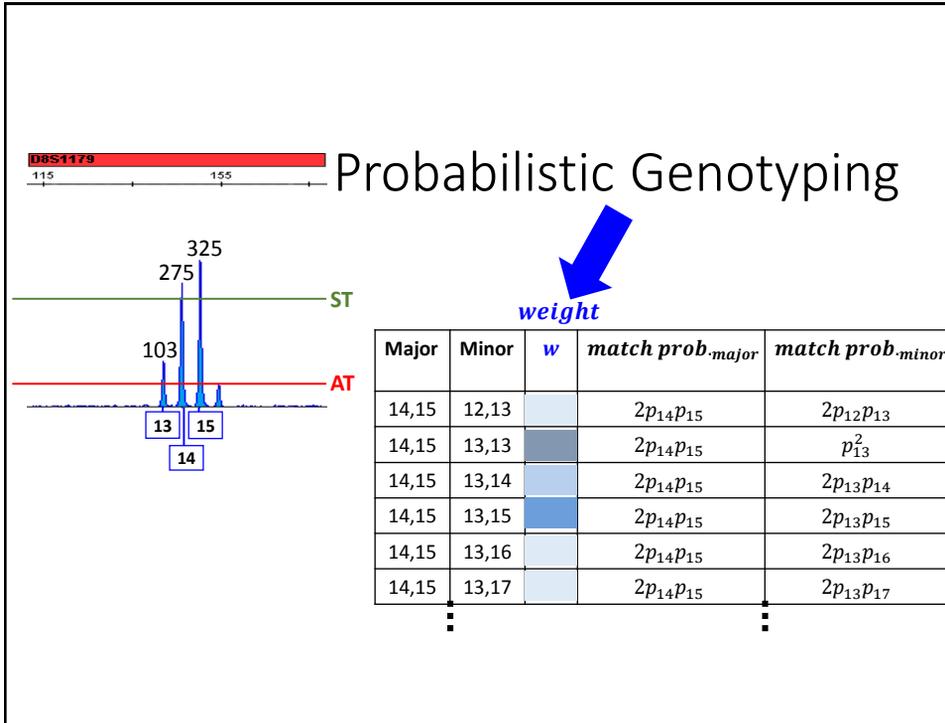
the probability of genotype set  $S_j$  given that we have observed genotypes  $G_K$  and that the contributors are as specified in  $H_p$  or  $H_d$

## Likelihood Ratio

$$LR = \frac{\sum_{j=1}^m f(G_{CS}|S_j) \Pr(S_j|G_K, H_p, I)}{\sum_{j=1}^m f(G_{CS}|S_j) \Pr(S_j|G_K, H_d, I)}$$

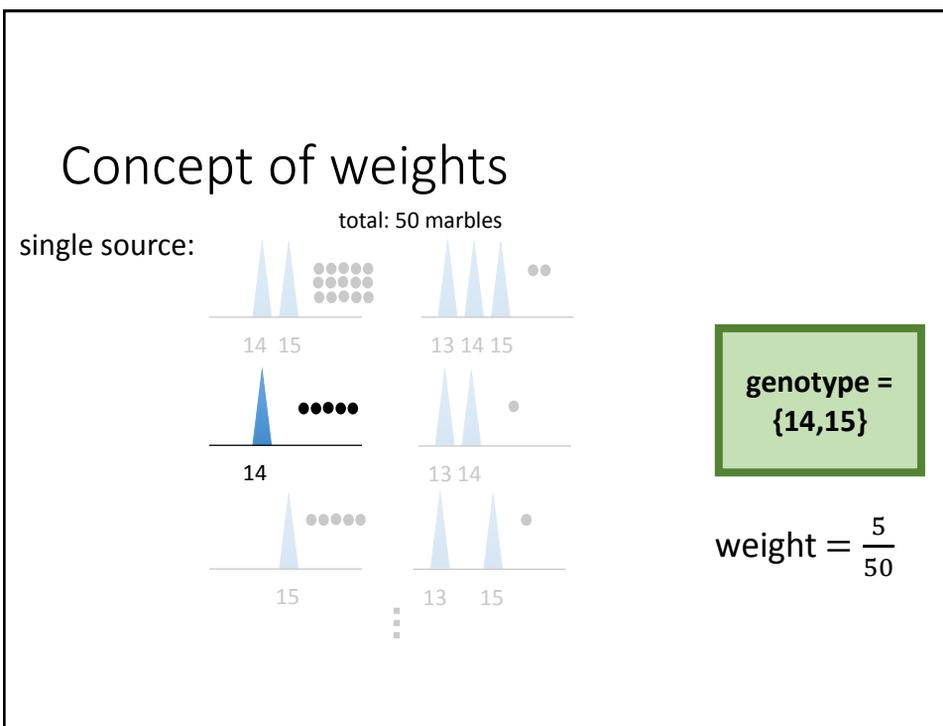
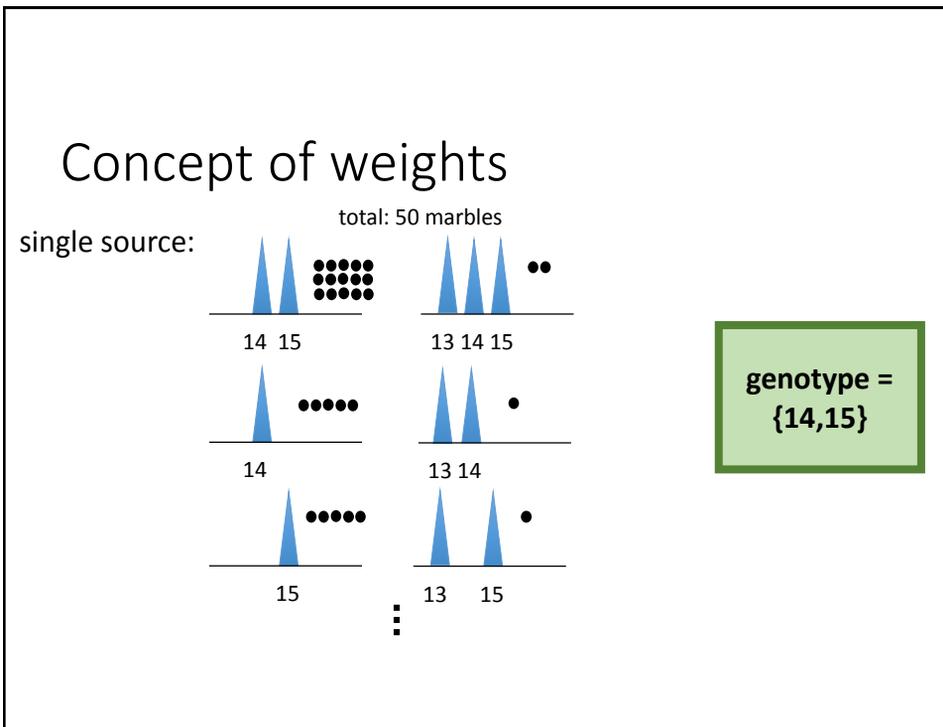
### weights

the probability of obtaining these DNA typing results for genotype set  $S_j$



## Concept of weights

the probability of obtaining these DNA typing results for a particular genotype set



## Concept of weights

how large or small a weight is depends on:

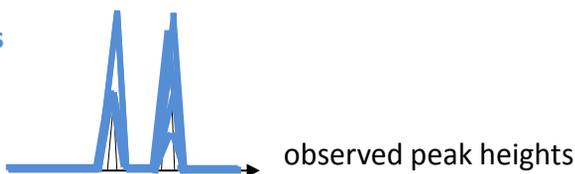
- possibility of allele drop-out
- possibility of allele drop-in
- amount of template DNA
- possibility of degradation
- possibility of oversized stutter peaks

⋮

## Continuous Model: short explanation

Simulations attempt to reproduce the observed peak heights by varying a set of parameters.

simulated peak heights



The better  $S_j$  and the set of parameters explain  $G_{CS}$ , the greater the probability density  $f(G_{CS}|S_j)$  assigned to that genotype set.

## Continuous model: longer explanation

- 1) Randomly choose a genotype set and values for all of the parameters that are necessary to model the expected peak heights.

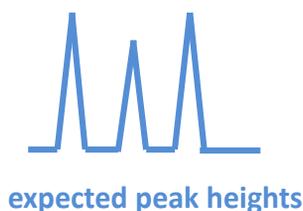
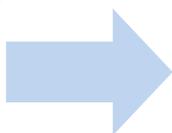
genotype set  $S_j$   
parameter 1  
parameter 2  
parameter 3  
⋮

**Markov Chain Monte Carlo (MCMC)  
sampling**

## Continuous model: longer explanation

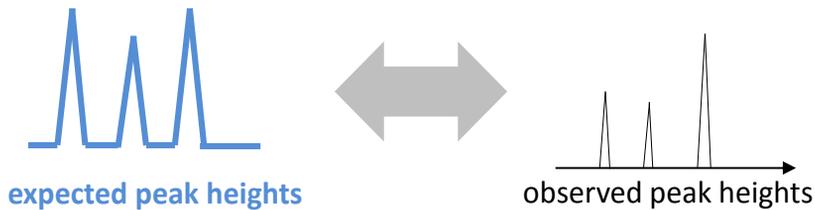
- 2) Model the expected peak heights given the chosen genotype set  $S_j$  and parameter values.

genotype set  $S_j$   
parameter 1  
parameter 2  
parameter 3  
⋮



## Continuous model: longer explanation

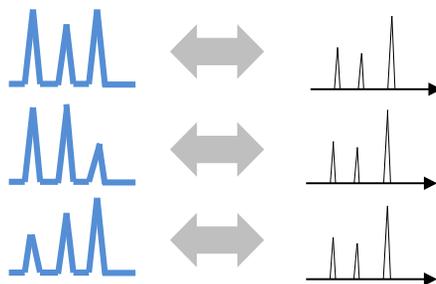
- 3) Compare the expected peak heights with the observed peak heights. How similar are they?



Assign a probability density for the observed peak heights given the expected peak heights.

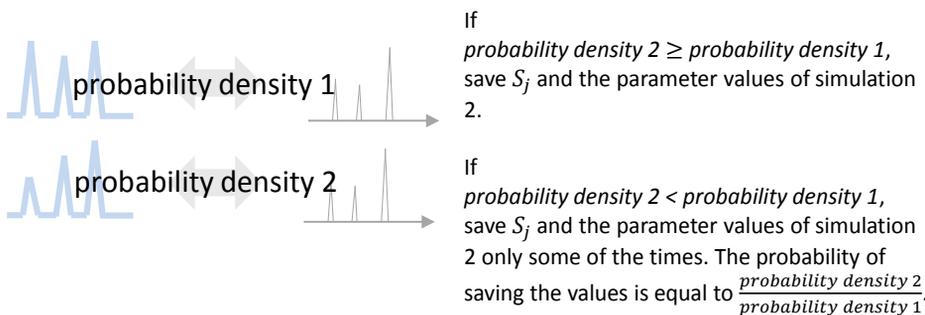
## Continuous model: longer explanation

- 4) Repeat steps 1) through 3) a large number of times. A large number of simulations are performed that randomly vary the genotype set and the parameter values.



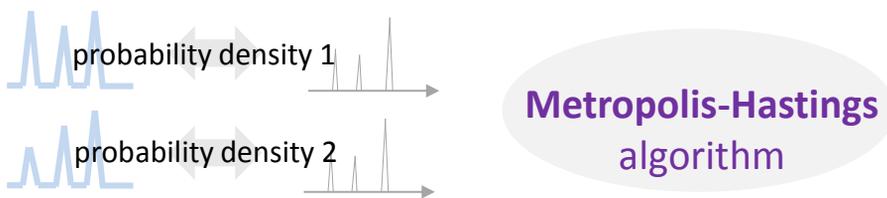
## Continuous model: longer explanation

- 3) Compare the expected peak heights with the observed peak heights. How similar are they?
- 4) Repeat steps 1) through 3) a large number of times. A large number of simulations are performed that randomly vary the genotype set and the parameter values.



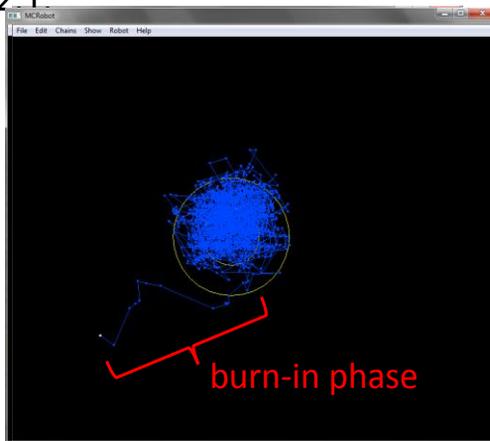
## Continuous model: longer explanation

- 3) Compare the expected peak heights with the observed peak heights. How similar are they?
- 4) Repeat steps 1) through 3) a large number of times. A large number of simulations are performed that randomly vary the genotype set and the parameter values.



## Continuous model: longer explanation

MCRobot-2.1:



after taking 300 samples

384

## Continuous model: longer explanation

5) Assign the weight  $f(G_{CS}|S_j)$  based on the saved simulation results.

saved simulation results:

Genotype set	Parameter 1	Parameter 2
<del>6,7 and 8,9</del>		
<del>6,7 and 8,8</del>		
<del>7,7 and 6,8</del>		
6,7 and 8,8		
6,7 and 8,8		
6,8 and 7,8		
6,7 and 8,8		
6,7 and 8,8		
6,8 and 7,8		

A red bracket on the left side of the table, labeled 'discard burn-in', encompasses the first three rows of the data.

## Continuous model: longer explanation

- 5) Assign the weight  $f(G_{CS}|S_j)$  based on the saved simulation results.

The weight  $f(G_{CS}|S_j)$  is assigned as the proportion of genotype sets  $S_j$  among the saved genotype sets.

Example:  $S_j = 6,7$  and  $8,8$

$$\rightarrow f(G_{CS}|S_j) = \frac{4}{6}$$

saved simulation results:		
Genotype set	Parameter 1	Parameter 2
<del>6,7 and 8,9</del>		
<del>6,7 and 8,8</del>		
<del>7,7 and 6,8</del>		
6,7 and 8,8		
6,7 and 8,8		
6,8 and 7,8		
6,7 and 8,8		
6,7 and 8,8		
6,8 and 7,8		

## Continuous model: longer explanation

- simulations
- 1) Randomly choose a genotype set  $S_j$  and values for all of the parameters that are necessary to model the expected peak heights.
  - 2) Model the expected peak heights given the chosen genotype set  $S_j$  and parameter values.
  - 3) Compare the expected peak heights with the observed peak heights. How similar are they?
  - 4) Repeat steps 1) through 3) a large number of times. A large number of simulations are performed that randomly vary the genotype set  $S_j$  and the parameter values.
  - 5) Assign the weight  $f(G_{CS}|S_j)$  based on the saved simulation results.

## Summary of main points

- The peak heights are a continuous variable.
- Mathematically speaking,  $f(G_{CS}|S_j)$  is the **probability density** for the observed peak heights given the genotype set  $S_j$ .
- The probability densities  $f(G_{CS}|S_j)$  are assigned based on the results of simulations that attempt to reproduce the observed peak heights by varying a set of parameters.